Clustering - create groupings/assignment of objects s.t
- similar to each other
- dissimilar to other clusters

USES
- Outlier/Anomaly detection (data cleaning | credit card fraud | spam filter)
- Filling gaps in Data (same marketing strategy for similar value | infer probable values for gaps)

Clustering Problem
- SIMILAR DATA POINTS IN SAME CLUSTER
- DISSIMILAR DATA POINTS IN DIFF CLUSTERS

- similar meaning?
- how to find clusters?
- whats a good cluster?

TYPES
- Partitional (each obj → 1 cluster)
- Heirarchical (set of nested clusters organized in a tree → (phylogenetic tree based on DNA/genome seq)
- Density-based (defined based on local density of points)
- Soft clustering (each point assigned to every cluster w/ a probability → (weight of species, → speak more in terms of likelihood)

PARTITIONAL

n data points, set to have K clusters up front

GOAL: partition n into K while maximizing similarity w/in cluster + dissim b/w clusters [intracluster]

intracluster dist :    $$\sum_{k=1}^{K} \sum_{x_i, x_j \, \in \, C_k} d(x_i, x_j)$$

pairwise distances per cluster

summed over all clusters.

NOT efficient so FIND CENTER OF MASS PER CLUSTER [CENTROID]

when d is euclidean, the centroid of m points is the average of the points. (use size of cluster |$C_k$| to weight the distances to the centroid.)

K-means
$$\sum_{i}^{K} \sum_{x \in C_i} \|x - \mu_i\|_2^2$$      FIND K points (means) that minimize the cost function.

NOTE: easy if k=1 or k=n

NOTE: if $x_i$ in greater than 2 dim ⇒ NP-hard.

K-means (Lloyd's Algo)
1. Randomly pick K centers ($\mu_1, \cdots, \mu_k$)
2. Assign each point in dataset to its closest center
3. Compute the new centers as means of each clusters.
4. Repeat 2,3 until convergence.

Proof by contradiction
    suppose it does not converge then
    1. minimum of cost func is only reached at limit
        CONTRADICTION: only finite data points, so we can't have infinite iterations
    2. cycle
        CONTRADICTION: suggests we increase cost func at some point.

CONCLUSION: ALWAYS converges
NOTE: will not always converge to optimal
    ↳ all depends on K points you start w/.
- Farthest First Traversal
    Choose centers as far from others → BUT does poorly w/ outliers.
- K-means++
    1. start w/ random center
    2. let D(x) be dist b/w x and centers selected so far, choose next center w/prob. proportional to $D(x)^a$
        a=0 → random        a=2 ⇒ k-means++        a=∞ → farthest first traversal
    generate a rand num 0 to N. and set intervals w/ prob of each point, and the value to be picked will be proportional
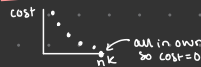        ↳ $\sum_{x \in D} D(x)^2$

K-means limitation
- splits large clusters
- dislikes nonglobular cluster shape
- dislikes varying density

How to choose right K?
1. iterate through different values of K (elbow method)
2. use empirical/domain-specific knowledge.
ELBOW METHOD



cost

n_k        all in own so cost=0

ELBOW METHOD MAXIMIZES COST AND NUMBER OF CLUSTERS.

NOTE: K-means is good for spherical gaussians

# K-mean variations

- K-medians (uses $L_1$ norm / manhattan dist)
- K-medoids (any dist func + centers must be in dataset)
- weighted k-means (each point has a different weight when computing mean) → good for outliers.